

WHAT IS CLAIMED IS:

Sub
A1

1. A method of categorizing an initial collection of documents, each document being represented by a string of characters, the method comprising the steps of:

identifying predefined characters in the string of characters from the documents in the initial collection of documents to form identified characters;
changing the identified characters in the documents in the initial collection of documents to form a preprocessed collection of documents;
constructing a number of categories from the preprocessed collection of documents; and
assigning each document in the preprocessed collection of documents to a category to form a hierarchy of categories of documents.

2. The method of claim 1 wherein the step of constructing a number of categories includes the steps of:

clearing a temporary category and selecting a seed document as a first document of the temporary category;
collecting documents from the preprocessed collection of documents that are similar to the seed document into the temporary category;
testing to determine if there are enough documents in the temporary category to merit construction of a new category;
constructing the new category and generating a heading for the new category if there are enough documents in the temporary category to merit construction;
assigning the seed document to a category reserved for documents not belonging to any specific category if there are not enough documents in the temporary category; and
marking the documents assigned to any category in the preprocessed collection of documents as processed.

09844040-042701

3. The method of claim 2 wherein the predefined characters include punctuation marks, and the changing step removes the punctuation marks from the string of characters.

4. The method of claim 2 wherein the predefined characters include upper-case characters, and the changing step replaces upper-case characters with lower-case characters.

5. The method of claim 2 wherein the predefined characters include non-root words, and the changing step replaces the non-root words with root words.

6. The method of claim 2 wherein the predefined characters include abbreviations, and the changing step replaces the abbreviations with original words.

7. The method of claim 2 wherein the predefined characters include articles, and the changing step removes the articles from the string of characters.

8. The method of claim 2 wherein the collecting step further includes the step of loading a character string from the seed document into a memory location to initialize the values of a number of category properties for the temporary category.

9. The method of claim 8 and further comprising the steps of:
determining if there are documents in the preprocessed collection of documents that have not been processed with respect to the temporary category;
if there are documents in the preprocessed collection of documents that have not been processed with respect to the temporary category, selecting a next document from the preprocessed collection of documents and measuring a

similarity with a similarity test between the selected document and a number of current category properties;

including the selected document in the temporary category if the selected document passes the similarity test;

updating the values of the number of category properties of the temporary category when the selected document is included; and

rejecting the selected document if the selected document fails the similarity test.

10. The method of claim 9 and further comprising the step of repeating the steps of claim 9 for all documents in preprocessed collection of documents.

11. The method of claim 2 wherein the collecting step further includes the step of collecting more similar documents from a number of existing categories.

12. The method of claim 11 and further comprising the steps of:
determining if there are more documents in a number of existing categories that have not been processed with respect to the temporary category;

if there are documents in the number of existing categories that have not been processed with respect to the temporary category, selecting a next document from the number of existing categories as a selected document and measuring a similarity with a similarity test between the selected document and a number of current category properties;

including the selected document in the temporary category if the selected document passes the similarity test; and

rejecting the selected document if the selected document fails the similarity test.

089984-0276267

13. The method of claim 12 and further comprising the step of repeating the steps of claim 12 for all documents in the number of existing categories.

14. The method of claim 8 wherein the category properties includes a string of characters selected from the group consisting of a longest common substring in the title, a longest common substring in the body; and a document type index measured as list of fractional numbers for each document type.

15. The method of claim 14 wherein a document type includes types selected from the group consisting of news article, technical documents, and poems.

16. The method of claim 2 and further comprising the steps of: making sub-categories if there are too many documents in a given category; and post-processing the number of categorized lists of documents.

17. The method of claim 16 wherein the categorized list of documents is post-preprocessed by the following steps:
merging two categories that each have a heading where there is too much overlap in the headings of the two categories; and
promoting sub-categories to an upper level in a hierarchy when there are not enough categories in the upper level.

18. The method of claim 2 wherein the seed document is a first document in the preprocessed collection of documents.

19. The method of claim 2 wherein the seed document is a document with a highest rank value among the documents not marked as processed in the preprocessed collection of documents.

20. The method of claim 2 wherein the temporary category is tested to determine if there are enough documents in the temporary category to merit construction of a new category by accumulating the weight of each document when each document can contribute uniform weight or different weight based on the rank value of each document with higher ranked document given more weight.

21. The method of claim 2 wherein the heading is a longest common substring in a title.

22. The method of claim 21 wherein the heading includes a number of longest common substrings.

23. The method of claim 1 and further comprising the steps of:
determining if an anchor-text character string is available for the documents in the initial collection of documents; and
attaching an anchor-text character string to the string of characters that represents the documents in the initial collection of documents when the anchor-text character string is available.

24. The method of claim 23 wherein the anchor-text character string is a text used most frequently by hypertext documents.

25. The method of claim 23 wherein the anchor-text character string is a text with a highest partial extrinsic rank value.

26. A method of categorizing an initial collection of documents, each document being represented by a string of characters, the method comprising the steps of:

constructing a number of categories from the initial collection of documents wherein a category is constructed by:

clearing a temporary category and selecting a seed document as a first document of a temporary category;

collecting documents from the initial collection of documents to the temporary category that are similar to the seed document;

testing to determine if there are enough documents in the temporary category to merit construction of a new category;

constructing the new category and generating a heading for the new category if there are enough documents in the temporary category to merit construction;

assigning the seed document to a category reserved for documents not belonging to any specific category if there are not enough documents in the temporary category; and

marking the documents assigned to any category in the initial collection of documents as processed; and

assigning each document in the initial collection of documents to a category to form a hierarchy of categories of documents.

27. The method of claim 26 wherein the collecting step further includes the step of loading a character string from the seed document into a memory location to initialize values of a number of category properties for the temporary category.

28. The method of claim 27 and further comprising the steps of:
determining if there are documents in the initial collection of documents that have not been marked as processed;

if there are documents in the initial collection of documents that have not been marked as processed, selecting a next document from the initial collection of documents and measuring a similarity with a similarity test between the selected document and a number of current category properties;

including the selected document in the temporary category if the selected document passes the similarity test; and

rejecting the selected document if the selected document fails the similarity test.

29. The method of claim 28 and further comprising the step of repeating the steps of claim 28 for all documents in initial collection of documents.

30. The method of claim 26 wherein the collecting step further includes the step of collecting more similar documents from a number of existing categories.

31. The method of claim 30 and further comprising the steps of:
determining if there are more documents in the number of existing categories that have not been processed with respect to the temporary category;
if there are documents in the number of existing categories that have not been processed with respect to the temporary category, selecting a next document from the number of existing categories and measuring a similarity with a similarity test between the selected document and a number of current category properties;

including the selected document in the temporary category if the selected document passes the similarity test; and

rejecting the selected document if the selected document fails the similarity test.

32. The method of claim 31 and further comprising the step of repeating the steps of claim 31 for all documents in number of existing categories.

33. The method of claim 1 wherein each document in the preprocessed collection of documents is assigned to one or more categories to form a hierarchy of categories.

34. The method of claim 26 wherein each document in the initial collection of documents is assigned to one or more categories to form a hierarchy of categories.

35. The method of claim 2 and further comprising the step of repeating the steps of claim 2 until all documents in the preprocessed collection of documents are marked as assigned to a category.

36. The method of claim 35 wherein the documents in the preprocessed collection of documents are initialized as unmarked before selecting a first seed document.

37. The method of claim 26 and further comprising the step of repeating the constructing steps of claim 26 until all documents in the initial collection of documents are marked as assigned to a category.

38. The method of claim 37 wherein the documents in the preprocessed collection of documents are initialized as unmarked before selecting a first seed document.

39. An apparatus that categorizes a collection of documents, each document being represented by a string of characters, the apparatus comprising:
means for identifying predefined characters in the string of characters from each document to form identified characters;
means for changing the identified characters in each document to form a preprocessed collection of documents;

means for constructing a number of categories from the preprocessed collection of documents; and

means for assigning each document in the preprocessed collection of documents to a category to form a number of categorized lists of documents.

089984-0276267